# An Automatic Finite-Sample Robustness Metric: Can Dropping a Little Data Make a Big Difference?

Ryan Giordano (`rgiordan@mit.edu`)[1]
January 2021

---

[1]With coauthors Rachael Meager (LSE) and Tamara Broderick (MIT)

## Dropping data: Motivation

You're a data analyst, and you've

- Gathered some exchangeable data,
- Cleaned up / removed outliers,
- Checked for correct specification, and
- Drawn a conclusion from your statistical analysis
  (e.g., based the sign / significance of some estimated parameter).

You're a data analyst, and you've

- Gathered some exchangeable data,
- Cleaned up / removed outliers,
- Checked for correct specification, and
- Drawn a conclusion from your statistical analysis
  (e.g., based the sign / significance of some estimated parameter).

**Well done!**

Would you be concerned if you could **reverse your conclusion** by removing a **small proportion** (say, 0.1%) of your data?

## Dropping data: Mexico Microcredit

Consider **?**, a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.
The variable "Beta" estimates the effect of microcredit in US dollars.

|          | Left out points | Beta (SE)     |
|----------|-----------------|---------------|
| Original | 0               | -4.55 (5.88)  |

Consider **?**, a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.
The variable "Beta" estimates the effect of microcredit in US dollars.

|             | Left out points | Beta (SE)      |
|-------------|-----------------|----------------|
| Original    | 0               | -4.55 (5.88)   |
| Change sign | 1               | 0.4 (3.19)     |

Consider **?**, a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.
The variable "Beta" estimates the effect of microcredit in US dollars.

|                     | Left out points | Beta (SE)     |
|---------------------|-----------------|---------------|
| Original            | 0               | -4.55 (5.88)  |
| Change sign         | 1               | 0.4 (3.19)    |
| Change significance | 14              | -10.96 (5.57) |

## Dropping data: Mexico Microcredit

Consider **?**, a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.
The variable "Beta" estimates the effect of microcredit in US dollars.

|                     | Left out points | Beta (SE)      |
|---------------------|-----------------|----------------|
| Original            | 0               | -4.55 (5.88)   |
| Change sign         | 1               | 0.4 (3.19)     |
| Change significance | 14              | -10.96 (5.57)  |
| Change both         | 15              | 7.03 (2.55)    |

## Dropping data: Mexico Microcredit

Consider **?**, a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.
The variable "Beta" estimates the effect of microcredit in US dollars.

|                     | Left out points | Beta (SE)      |
|---------------------|-----------------|----------------|
| Original            | 0               | -4.55 (5.88)   |
| Change sign         | 1               | 0.4 (3.19)     |
| Change significance | 14              | -10.96 (5.57)  |
| Change both         | 15              | 7.03 (2.55)    |

By removing very few data points ($15/16560 \approx 0.1\%$), we can reverse the qualitative conclusions of the original study!

Consider **?**, a randomized controlled trial study of the efficacy of microcredit in Mexico based on 16,560 data points.
The variable "Beta" estimates the effect of microcredit in US dollars.

|                     | Left out points | Beta (SE)      |
| ------------------- | --------------- | -------------- |
| Original            | 0               | -4.55 (5.88)   |
| Change sign         | 1               | 0.4 (3.19)     |
| Change significance | 14              | -10.96 (5.57)  |
| Change both         | 15              | 7.03 (2.55)    |

By removing very few data points ($15/16560 \approx 0.1\%$), we can reverse the qualitative conclusions of the original study!
**Question:** Is the reported interval $-4.55 \pm (5.88)$ a reasonable description of the uncertainty in the estimated efficacy of microcredit?

## Dropping data: Motivation

Would you be concerned if you could **reverse your conclusion** by
removing a **small proportion** (say, 0.1%) of your data?

## Dropping data: Motivation

Would you be concerned if you could **reverse your conclusion** by
removing a **small proportion** (say, 0.1%) of your data?
**Not always!**

## Dropping data: Motivation

Would you be concerned if you could **reverse your conclusion** by
removing a **small proportion** (say, 0.1%) of your data?
**Not always!**
**...but sometimes, surely yes.**
For example, often in economics:

- Small fractions of data are missing not-at-random,
- Policy population is different from analyzed population,
- We report a convenient summary (e.g. mean) of a complex effect,
- Models are stylized proxies of reality.

## Objective

Estimate the effect of leaving out $\lfloor \alpha N \rfloor$ datapoints, where $\alpha$ is small.

## Objective

Estimate the effect of leaving out $\lfloor \alpha N \rfloor$ datapoints, where $\alpha$ is small.

**Question 1: How do we find influential datapoints?**
The number of subsets $\binom{N}{\lfloor \alpha N \rfloor}$ can be very large even when $\alpha$ is very small.
In the MX microcredit study, $\binom{16560}{15} \approx 1.4 \cdot 10^{51}$ sets to check for $\alpha = 0.0009$.
We provide a fast, automatic approximation based on the **influence function**.

## Objective

Estimate the effect of leaving out $\lfloor \alpha N \rfloor$ datapoints, where $\alpha$ is small.

**Question 1: How do we find influential datapoints?**

**Question 2: What makes an estimator non-robust?**

## Objective

Estimate the effect of leaving out $\lfloor \alpha N \rfloor$ datapoints, where $\alpha$ is small.

**Question 1: How do we find influential datapoints?**

**Question 2: What makes an estimator non-robust?**

Non-robustness to removal of $\lfloor \alpha N \rfloor$ points is:

- Not (necessarily) caused by misspecification.
- Not (necessarily) caused by outliers.
- Not captured by standard errors.
- Not mitigated by large $N$.
- Primarily determined by the **signal to noise** ratio
    ... in a sense which we will define.

## Objective

Estimate the effect of leaving out $\lfloor \alpha N \rfloor$ datapoints, where $\alpha$ is small.

**Question 1: How do we find influential datapoints?**

**Question 2: What makes an estimator non-robust?**

**Question 3: When is our approximation accurate?**

## Objective

Estimate the effect of leaving out $\lfloor \alpha N \rfloor$ datapoints, where $\alpha$ is small.

**Question 1: How do we find influential datapoints?**

**Question 2: What makes an estimator non-robust?**

**Question 3: When is our approximation accurate?**

- We provide deterministic error bounds for small $\alpha$.
- We show the accuracy in simple experiments.
- We show the accuracy in a number of real-world experiments.

## Objective

Estimate the effect of leaving out $\lfloor \alpha N \rfloor$ datapoints, where $\alpha$ is small.

**Question 1: How do we find influential datapoints?**

**Question 2: What makes an estimator non-robust?**

**Question 3: When is our approximation accurate?**

**Conclusion: Related work and future directions**

**Question 1:**
**How do we find influential datapoints?**

**Question 2:**
**What makes an estimator non-robust?**

**Question 3:**
**When is our approximation accurate?**

**Conclusion:**
**Related work and future directions**